

4



Module Leaders



Jan Smeddinck jan.smeddinck@ncl.ac.uk https://openlab.ncl.ac.uk/people/jan-smeddick/



Yu Guan yu.guan@ncl.ac.uk https://openlab.ncl.ac.uk/people/yu-guan/





5

Motivation for this module...

- Most AI/ML courses are disconnected from real-world problems...
- Most AI/ML courses consider "user-interfaces" or human impact as an afterthought; and focus narrowly on algorithms...
- ... why is this a problem?

Gender	Gender Classifier	Darker Male	Darker Female	Lighter Male	Lighter Female	Largest Gap
snades	Microsoft	94.0%	79.2%	100%	98.3%	20.8%
Buolamwini, J., &	FACE**	99.3%	65.5%	99.2%	94.0%	33.8%
Gebru, I. (2018). Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. 15.	IBM	88.0%	65.3%	99.7%	92.9%	34.4%





Further reading...







Acknowledgements

- This module builds on many existing sources and benefits from the generous contributions by a number of individuals...
- Key existing module: <u>https://haiicmu.github.io</u> by <u>Haiyi Zhu</u> and <u>Steven Wu</u>; formerly by <u>Chinmay Kulkarni and Mary</u> <u>Beth Kery</u>
- Individual contributions (lecture materials or sessions) by Alex Bowyer, Robert Porzel, Viana (Nijia) Zhang, and more.

🔭 Interdisciplinary College

Neuroscience • Neuroinformatics • Cognitive Science • Artificial Intelligence • and more

Annual interdisciplinary spring school virtual (affordable!) version in 2021 (staring March 12). Theme: Connected in Cyberspace



https://interdisciplinary-college.org/







Course Outline: Week 03

- Topic 01 Explainable, Interpretable & Relatable AI
- Topic 02 AI Ethics
- Topic 03 Humans-in-the-Loop
- Topic 04 Recommender Systems
- Topic 05 Conversational Interfaces
- Topic 06 Al Agents & Robots
- Topic 07 Human-Al Integration
- Topic 08 Creative Al
- Topic 09 Summary & Outlook



Deep Fake Sandbox

Fun video introduction link: https://www.youtube.com/ watch?v=mUfJOQKdtAk



Siarohin, A., Lathuilière, S., Tulyakov, S., Ricci, E., & Sebe, N. (2019). First Order Motion Model for Image Animation. Advances in Neural Information Processing Systems, 32, 7137–7147.

https://colab.research.google.com/github/AliaksandrSiarohin/fir st-order-model/blob/master/demo.ipynb https://aliaksandrsiarohin.github.io/first-order-model-website/

21



Learning Goals

Learn / get a refresher on the foundational terms behind...

- Artificial Intelligence
- Machine Learning
- (Non-Learning) Adaptive Systems
- ... and their most general components / elements

24









Machine Learning

A "machine" that is able to improve based on past experience without explicit human programming on how to improve each time.

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves with experience E."

-- Tom M. Mitchell

n









32































-	Classification - Is this a dog or a cat? Regression - How warm will it be tomorrow?	Supervised learning
-	Clustering - Here are a few articles. Organize them into topics Representation - What is the best way to represent words? (e.g. so representation "autumn" and "fall" are similar) - What is the color of happiness?	Unsupervised learning of
-	"Dear algorithm: figure it out, and if you're right, you get - Drive a car - Recommend movies to people	a reward" Reinforcemen learning















Prominent Algorithms Summary

- NN
- Regression analysis
- K-means
- Genetic algorithms
- Bayesian networks
- Decision trees
- ... more (and more details) in week 2!

57



<section-header>
 Applications: Natural Language Processing
 Applications: Natural Language Processing
 Rich history ... may or may not make use of Al/ML
 Understanding) or NLG (generation)
 Machine Translation
 Sentiment Analysis
 Machine Translation
 Lest Summarization
 Senticationi
 Speech Recognition



60







Applications: Computer Vision



This Photo by Alex Chitu is licensed under CCBY

Dean, J., Corrado, G. S., Monga, R., Chen, K., Devin, M., Le, Q. V., Mao, M. Z., Ranzato, M., Senior, A., Tucker, P., Yang, K., & Ng, A. Y. (2012). Large Scale Distributed Deep Networks. *NIPS*.



This Photo by Alex Chitu is licensed under CC BY

James Charles, Stefano Bucciarelli and Roberto Cipolla. 'Real-time screen reading: reducing domain shift for oneshot learning.' Paper presented at the British Machine Vision Conference. 2020

64





Ubiquitous Computing \rightarrow AI

- Dynamic / learning / adaptive systems becoming ubiquitous (esp. online services / mobile devices)
- AI/ML also becoming ubiquitous
- Commodity
- Still centralised (esp. amongst FAANG)
- Computing as "tools" to "services" and even "servants"
- "There's an app for that" \rightarrow "There's an agent for that"





Evaluating Learning Systems

Generally interested in:

- How often is the prediction wrong?
- How is the prediction wrong?
- What is the cost of wrong predictions?
- $\circ\,$ How does the cost vary by the type of prediction that was wrong?
- How can we minimize cost?

• HCAI: misunderstandings? frustration?regret?

70













via https://haiicmu.github.io

According to a preliminary report released by the National Transportation Safety Board last week, Uber's system detected pedestrian Elaine Herzberg six seconds before striking and killing her. It identified her as an unknown object, then a vehicle, then finally a bicycle. (She was pushing a bike, so close enough.) About a second before the crash, the system determined it needed to slam on the brakes. But Uber hadn't set up its system to act on that decision, the NTSB explained in the report. The engineers prevented their car from making that call on its own "to reduce the potential for erratic vehicle behavior." (The company relied on the car's human operator to avoid crashes, which is a whole separate problem.)

> Uber's engineers decided not to let the car auto-brake because they were worried the system would overreact to things that were unimportant or not there at all. They were, in other words, very worried about false positives.

https://www.wired.com/story/self-driving-carsuber-crash-false-positive-negative/ 75

ARN MORE









Silent Errors? UX?

- Not all errors are visible to the computer
- How do you measure when your device can't hear "OK, Google"
- More in later sessions...

Biases & Fairness

- Bias as in prejudice in favor or against a person, group, or thing that is considered to be unfair.
- Frequently see bias in the classifications and predictions
- Often due to how data is sampled / collected
- Can also be introduced in processing
- Many algorithms/methods now in development to analyse / check for biases / fairness
- Proper judgment may require humans-in-the-loop (more in week 3)

81

	0			1
	Vendor name	Funding	# of employees	Location
	8 and Above	141	1-10	WA, USA
	ActiView	\$6.5M	11-50	Israel
	Applied	£2M	11-50	UK
	Assessment Innovation	\$1.3M	1-10	NY, USA
1	Good&Co	\$10.3M	51-100	CA, USA
1	Harver	\$14M	51-100	NY, USA
1	HireVue	\$93M	251-500	UT, USA
1	impress.ai	\$1.4M	11-50	Singapore
1	Knockri		11-50	Canada
	Koru	\$15.6M	11-50	WA, USA
	LaunchPad Recruits	£2M	11-50	UK
	myInterview	\$1.4M	1-10	Australia
	Plum.io	\$1.9M	11-50	Canada
	PredictiveHire	A84,3M	11-50	Australia
	pymetrics	\$56.6M	51-100	NY, USA
1	Scoutible	\$6.5M	1-10	CA, USA
	Teamscope	6800K	1-10	Estonia
1	ThriveMap	LISIK	1-10	CA LEA
	robs	\$1M	11-50	CA, USA

BLOW B.08 s z 4 s s to Chardy Report Mar Manded ¹ Notice the space of the spac
CANDIDATE BLUEPRINT
Cypenness Endustan Warmees Kindness Reactions
Advecting
Otherst
Sente
Penouold
Field
Charge agent

82

81







Human-Al
htteraction &
htteraction &
htteraction &
hteraction &
hteracti





Don't Promise What You Can't Keep

"Machines will be capable, within twenty years, of doing any work a man can do." - Herbert Simon 1965

"In from three to eight years we will have a machine with the general intelligence of an average human being." - Marvin Minsky 1970

89









Robots #2 (1942)

1942 Isaac Asimov's Robot Laws

- 1. A robot may not injure a human being or, through inaction, allow a human being to come to harm.
- 2. A robot must obey the orders given it by human beings except where such orders would conflict with the First Law.
- 3. A robot must protect its own existence as long as such protection does not conflict with the First or Second Laws.
- More in later sessions...





<section-header><section-header><list-item><list-item><list-item>

98

<section-header><text><list-item><list-item><list-item><text>

Skinner teaching machine (1960)

• Roots of modern "e-learning" ...



https://www.youtube.com /watch?v=CFYruzWeFwQ

00

Engelbart (1968) (and McLuhan)

- "Mother of all demos" ... early HCI ...
- "A Research Centre for Augmenting Human Intellect"



https://www.youtube.com watch?v=Xptc6f3Daoo

101



102

1969 "Perceptron" (destroyed)

- Marvin Minsky and Seymour Papert publish a book "Perceptrons"
- Burn piece of Perceptron approach in favour of rule approach
- Perceptron cannot handle XOR operator
 - However, multi-layer perceptrons CAN (and this was known at the time)
- Shuts down funding for neural networks
- Rosenblatt soon dies, never to see neural nets revindicated

Symbolic Al

- Symbolism: formal logic systems can represent intelligent action
- Assumption: "Every aspect of learning or any other feature of intelligence can be so precisely described that a machine can be made to simulate it."
- Later to be enshrined in Lisp-Machines

```
(RULE 5
(IF (PCS-SCS HEAT TRANSFER INADEQUATE)
       (LOW FEEDWATER FLON))
(THEN (ACCIDENT IS LOSS OF FEEDWATER)))
      (IF (SG INVENTORY INADEQUATE)
(LOW FEEDWATER FLOW))
(THEM (ACCIDENT IS LOSS OF FEEDWATER)))
```

```
(RULE 7
```

```
(IF (PCS INTEGRITY CHALLENGED)
(CONTAINMENT INTEGRITY CHALLENGED))
(THEN (ACCIDENT IS LOCA)))
```

```
(RULE &
(1# (PCS INTECRITY CHALLENCED)
       (SG IFVEL INCREASING))
(THEN (ACCIDENT IS STEAM GENERATOR TURE
       RUPTURE )))
```

```
CRULE S
       (IF (SC INVENTORY INADEQUATE)
       (HIGH STEAK FLOW))
THEN (ACCIDENT IS STEAN LINE BREAK))))
```

Figure 2. Event-oriented IF-IEEN rules.

Symbolic Al

- Newell & Simon's General Problem Solver can solve math proofs by searching a logic space
- Advances in natural language processing based on rules
 how words relate
- Advances in computer vision based on image transforms
- Advances in robotics based on rules and search in simplified settings
- BUT: Strong limitations in scalability!

105

1st Al Winter (~1974)

Funding is lost due to unmet promises:

- Lighthill Report 1973 shuts down funding in UK
 - James Lighthill (1973): "Artificial Intelligence: A General Survey" in Artificial Intelligence: a paper symposium, Science Research Council
- DARPA switches to "mission-oriented direct research, rather than basic undirected research"
- Dreyfus at MIT argues lots of human reasoning is not based on logic rules, involving instinct and unconscious reasoning

 (No AI researcher will eat lunch with Dreyfus for the next decade)
- Sussman: "using precise language to describe essentially imprecise concepts doesn't make them any more precise."

106

Limitations with Logic rules: Combinatorial Explosion

- Representation of knowledge is hard
- All the possibilities (mutability vs. immutability of the sign ... semioticians such as de Saussure, Charles Peirce et al.)
- Combinatorial explosion and search space

• Common sense

 Robot motion planning fails: <u>https://www.youtube.com/</u> <u>watch?v=g0TaYhjpOfo</u>



1980's Al Spring: From General Intelligence to Specific Applications Expert systems: while rules can't do everything they

- Expert systems: while rules can't do everything the are useful for capturing what experts do
- 1989 Deep Thought at CMU can play human-level chess

Expert Systems (great prototypes, but riddled by scalability issues)





Again: hype gets too high and funding is lost!
 → 2nd Al winter (into 90s / 2000s ... depending on
 perspective)

1990s - 2000s Al "undercover"

- Many advances in computer science, probability theory, statistical learning, etc.
- Video games boom births GPUs...
- "Computer scientists and software engineers avoided the term artificial intelligence for fear of being viewed as wild-eyed dreamers." - NYT 2005
- Machine Learning as a field distances itself from AI (temporarily!)
- Online search and advertising push for personalisation / recommendations / user models
- General public "goes online" ... lots of data becomes available
- Structured approaches in "linked data" / ontologies

109



110

109





History of Human-Al Interaction

2 Reads:

- Shneiderman, B., & Maes, P. (1997). Direct manipulation vs. Interface agents. *Interactions*, 4(6), 42–61. https://doi.org/10.1145/267505.267514
- Horvitz, E. (1999). Principles of mixed-initiative user interfaces. Proceedings of the SIGCHI Conference on Human Factors in Computing Systems the CHI Is the Limit -CHI '99, 159–166. <u>https://doi.org/10.1145/302979.303030</u>

• Akin to Shneiderman's Golden Rules / Nielsen's Heuristics in UB

History of Human-Al Interaction

Horvitz, E. (1999). Principles of mixed-initiative user interfaces. Proceedings of the SIGCHI Conference on Human Factors in Computing Systems the CHI Is the Limit - CHI '99, 159–166. <u>https://doi.org/10.1145/302979.303030</u>

• Akin to Shneiderman's Golden Rules / Nielsen's Heuristics in UB

Critical factors for the effective integration of automated services with direct manipulation interfaces include:

- 1. Developing significant value-added automation
- 2. Considering uncertainty about a user's goals
- 3. Considering the status of a user's attention in the timing of services

4. ...

113

114

Machine Learning Boom (2010ish)

- By 2010s reached serious scale
- Moore's law + data, data, data
- Machine learning now immensely powerful (perception)
- Neural networks now practical, many earlier inventions in neural networks (e.g. back propagation from 1980s) are "re-discovered"
- Deep learning gives big performance leaps in almost all application areas overnight
- Massive funding is back (now combined research + "real-world")
- ML "re-unites" with AI ... starts to become pervasive / ubiquitous
- HCI has since gone from UB → UX ... with widespread use, neglect of user-centred approaches in learning systems becomes obvious







Standard Open Data Sources

- Traditional methods / algorithms development in Al frequently uses existing data sets (e.g. IRIS, MINST, ImageNet, etc.) ... i.e. published research data
 - Readily available; often well formatted / relatively complete
 - Good for benchmarking
- Cost / complexity of acquiring / processing / managing data often overlooked
- Applied AI/ML often needs real-world relevant data sets and builds on "user data"...
 - Can be difficult to get by, unorganised, sparse, legal implications
 - Difficult to share, hence difficult to compare



Justification Practices

Open Lab

From ongoing work by Kieran Cutting (supervisee):

- Design under austerity
- Work with young people in state foster care at transition to adulthood
- HCI now embraces postcolonialism, feminism, and anarchism
- Little associated methodological innovation
- Researchers to rely upon old methods in new contexts
- Justification practices: grounded theory how work, relationships, and experiences of care have changed as a result of austerity
- Can subsume well-intended technology design and use

K.L.Cutting2@newcastle.ac.uk



121

122

121



Sensing Devices (Wearables & Other)

Search & Stats, Advertising, Social Media

• Historically as enabled by

data" initiatives...

advent of (broad) internet use

• Also public data (beyond pure

research data) under "open

The Intersection of

Social Media and Big Data

From http://blogs.zdnet.com/Hinchcliffe

This Photo by Dion Hinchcliffe is licensed under CC BY-SA

• Rapid development in sensor technologies ... e.g. accelerometers in 1950s ~ \$100k+... today: cents

This Photo by Laura James is licensed under CC By













Human-Data Interaction (HDI)

- Elmqvist (2011): the "human manipulation, analysis, and sensemaking of large, unstructured, and complex datasets."
- HDI is about federating disparate personal data sources and enabling user control over the use of "my data" (McAuley, Mortier, and Goulding 2011)
- Relatively young area, so terminology and definitions still emerging / changing

Elmqvist, N., 2011. Embodied Human-Data Interaction. ACM CHI 2011 Work. "Embodied Interact. Theory Pract. HCI 104–107.

Human-Data Interaction (HDI)

Mortier et al. (2014): focus on personal and open data:

- <u>Legibility</u>: process of understanding and making data and analytics algorithms both transparent and comprehensible to people
- <u>Agency</u>: power of handling data; capacity to control, inform, and correct data and inferences
- <u>Negotiability</u>: dynamic relationships that emerge regarding data; understanding and attitudes change over time; reevaluate

Mortier, R., Haddadi, H., Henderson, T., McAuley, D., Crowcroft, J., 2014. Human-Data Interaction: The Human Face of the Data-Driven Society. SSRN Electron. J. https://doi.org/10.2139/ssrn.2508051

133

HDI Research Topics & Challenges

- Personal data, data ownership and consent
- Embodied interactions
- Data visualization, mining, and analytics [actually 3 topics]
- HDI for specific domains (health informatics, smart cities, ..) Challenges (selected):
- Users' engagement through participation in design processes
- Models and value of data ownership, social / cultural aspects
- Transcending human and machine limitations in data analysis
- Data influence in decision-making process
- Systemic view of the complete data life cycle

Victorelli, E. Z., Dos Reis, J. C., Hornung, H., & Prado, A. B. (2020). Understanding human-data interaction: Literature review and recommendations for design. International Journal of Human-Computer Studies, 134, 13–32. https://doi.org/10.1016/j.ijhcs.2019.09.004

134



<text><text><text><text><figure>

Key Considerations

- Transparency and audit
 - What audit trails and information are to be provided to support this?
- Privacy and control
 - How can the resulting audit data be used to enable interaction around control of access to and processing of data?
- Analytics and commerce
 - How can the analysis algorithms that are used be made transparent to users (often while retaining protected commercial status)?
- Data to knowledge
 - How can the vast amount of data be used to benefit the individuals and let the society exploit the wealth of information offered by shared data?

via https://haiicmu.github.io











Post-Mortem Privacy Paradox

- work with users of password managers to explore views on the sharing, security and privacy of common digital assets
- when facing loss of control (e.g. care towards end of life and death)
- users recognise value in planning for their digital legacy
- yet: avoid actively doing so
- tension between the use of recommended security tools during life and facilitating appropriate post-mortem access to chosen assets

Jack Holt, James Nicholson, & Jan Smeddinck. (2021, accepted). From personal data to digital legacy: Exploring conflicts in the sharing, security and privacy of post-mortem data. WWW '21: Proceedings of The Web Conference 2021. WWW '21: The Web Conference 2021, Lyubjana, Stovenia

143

Data Rights Initiatives & Law (e.g. GDPR)

- rules relating to the protection of natural persons with regard to the processing of personal data and rules relating to the free movement of personal data
- protects fundamental rights and freedoms of natural persons and in particular their right to the protection of personal data
- E.g. reason to store / process data, consent, simple terms, clear processing guidelines, consider time requirements, data subjects / controllers / processors, EU (or equivalency) vs. non-EU, ...
Data Rights Initiatives & Law (e.g. GDPR)

- Right to be forgotten
 - ... right to have data records pertaining to individuals removed (e.g. from search engines)
 - ... but what about impact on group models / derived data?
- Transparency of use → some level of "right to explainability" (more on that in week 3)
- Generally also legal protections "against biased models" based on discrimination protection
- May appear "hindering" from engineering POV, but generally helpful towards explainable and fair Al ...

145



Open Lab

See session video on

146

Personal Data Use

From ongoing work by Alex Bowyer (supervisee):

- Data is used in care, but it disempowers
- Ppl already disempowered, staff/gov/others want more data & more linkage
- Focus on staff needs, staff efficiency
 not looking at relationships
- HDI, agency
 - data consolidation will make it worse need shared HDI

 \rightarrow Digital Civics, shared stakeholder needs, find holistic way forward also GDPR data explorer

A.Bowyer2@newcastle.ac.ul



145





Human-Data Interaction Summary

An emerging (sub-)field: HDI is about ...

- federating disparate personal data sources and enabling <u>user control</u> over the use of "my data" (McAuley, Mortier & Goulding 2011)
- human manipulation, analysis, and sense-making of large, unstructured, and complex datasets (Elmqvist 2011)
- processes of <u>collaboration</u> with data and the development of communication tools that enable interaction (Kee et al. 2012)
- providing access and understandings of data that is about individuals and how it affects them (Mashhadi, Kawsar & Acer 2014)

https://www.interaction-design.org/literature/book/the-encyclopedia-of-humancomputer-interaction-2nd-ed/human-data-interaction Human-Data Interaction Summary actions include feeding back inferences as input data for subsequent analysis analytics inferences, often opaque to users, are generated and used to drive actions actions based on our data and that of others affect our subsequent behaviour we lack legibility we lack negotiability we lack agency means to manage "our" data and access means to navigate data's social means to inspect and reflect on "our" data, to understand its collection and processing to it, enabling us to act effectively in these systems as we see fit aspects in collaboration with others and their policies Author/Copyright holder: Richard Mortier. Copyright terms and licence: CC BY-NC-ND 150

150



































Visualisation Empowerment

- Democratising data science → enables civic knowledge and awareness, but also argumentation; this often uses visualisations
- See e.g. work on visualisation empowerment (and VR Data Vis & Ix) by Benjamin Bach (University of Edinburgh)

Further reading:

- Data Visualisation Handbook (Koponen & Hilden)
 https://datavizhandbook.info/
- Savvy for more?: <u>https://datavis-online.github.io/</u>

169

Visualization for developing AI/ML systems

- Play with...
- <u>https://research.google.com/bigpicture/attacking-discrimination-in-ml/</u>
- What-if tool: <u>https://pair-code.github.io/what-if-tool/</u>
- What do you think of these visualizations / tools?
- Where you able to gain insights? Why? Why not?

170

169





Human-Al Interaction

An emerging term...

- Human-Al Interaction (HAI/HAII)
- Human-Centred/Centric AI (HCAI)
- Al Interaction Design (AIxD)
- AI User Experience (AIUX)

... differences in emphasis and application focus ...



Put-that-there (1980)

- How could this be helpful (augment)?
- How could this fail?
- Is it efficient?
- Is it a good UX?

Bolt, R. A. (1980). \"Put-that-there\": Voice and gesture at the graphics interface. Proceedings of the 7th Annual Conference on Computer Graphics and Interactive Techniques, 262–270. 174

174



Al Agents in 2020 (preview)





Personalisation & Customization

Concept	Definition	
Adaptability	The fact that a system is not fixed, but can be changed (to the needs of users, to changing environmental contexts, etc.; changes are usually understood to be performed manually).	
Customization	The act of changing a system to the needs of a user group or individual user (manually or automatically; may can be done by the group itself or by the user him- or herself, but may also be done by third parties; often related to the appearance or content of the given system).	
Personalization	The act of changing a system to the needs of a specific indi- vidual user (often automatic but does not have to be, i.e., can be understood as a specific form of customization with a focus on individuality; personalization is also often related to appearance or content).	Note: can be heuristic, ML based, etc.
Adaptivity	The fact that a system is not fixed, but dynamically changes over time (to adjust to the needs of users or an individual user, or to adjust to changing environmental contexts, etc.; typically happens automatically; often related to settings and parameters present in the given system).	Smeddinck, J. (2017). Human- Computer Interaction with Adaptable & Adaptive Motion- based Games for Health [University of Bremen]. SUUB Bremen. https://arxiv.org/abs/2012.03309

178



Adaptive Interfaces: Where is Audacity !?!

















28/12/2021







Learning from Interactions for Adaptive Systems

Heuristic adaptive systems usually "reactive" ... to more effectively prevent going in "wrong directions" need to be "predictive" \rightarrow Al/ML

- Build knowledge about the user, the interface, the agent and the world in terms of properties and capabilities...
- Learn to interpret the user behavior...
- Predict their needs and desires...
- Reason on the outcomes of different potential actions...
- Aka: contextual computing with context models and feedback...







UB/UX vs Uncertainty & Unpredictability

• UB/UX "laws" vs. "Al" systems

Links to direct manipulation...

- Continuous Representation of Object of Interest (feedback)
- Rapid, incremental, and reversible interactions
- Physical Actions

Potential & Risks

- Relinquishing control to an AI/ML agent can be helpful, but can be much harder to correct or understand if things go wrong
- "Unpredictability" can be a positive thing in one kind of experience, and a terrible thing in another...





198

197

What was the error?: Severe Failure Any's earlier version Xiaolce ran on China's most widespread instant messaging app Wechat ... without any major ethical incidents What makes Twitter a different environment? Any had no moral agency. To her, words like Hitler or Holocaust are not different from words like chair or Oklahoma 2018 version used black-listing and "moral judgments"... better? Cencorship? Nuance? Same still true for very large language models (ala GPT-3)?



Other Errors: Poor model performance

- Usually solvable by acquiring more training data for the situations the model is weakest at
- Data is expensive to collect and prepare, and your company or organization has limited resources.
 Prioritizing what specific data to collect is essential
- Designers can use rule or non-ML based fallbacks to still deliver the user some value when model performance isn't good enough for some cases
- Mixed models (Wizard-of-Oz) often used in business world ... in particular start-ups...

201

Other Errors: Low confidence or *false* High confidence in a prediction

- Low confidence predictions can mean that the model has lower performance, or the phenomena itself is just... less predictable
- Communicating with the user or providing good non-AI/ML fallbacks is key
 - Think levels of control (can be dynamic)
- High confidence (when the model is really wrong) is worse
 - Unkown unkown errors
- Need to give the user some error correction or feedback method to deal when this happens

202

201

Other Errors: Relevance/appropriateness errors

- Airbnb suggesting 'fun local activities' when you're traveling for a funeral
- Exercise app suggesting 'time to get up and walk!' when you're seated on a long flight
- Amazon suggesting products that you are allergic to or can't eat

Respectful Diet Information and Decision Support

From ongoing work with Remco Benthem de Grave (supervisee):

- Current diets are unsustainable
 environmental + health perspective
- Personalized support with digital tech.
- Available systems often coercive (overly persuasive)

→ Focus: supporting (non-coercively) existing diet values and intentions through personalisation and minimizing effort of making sustainable choices

r.benthemdegrave2@newcastle.ac.uk



Silent Errors? UX?

- Not all errors are visible to the computer
- How do you measure when your device can't hear "OK, Google"
- How does your fitness tracker know it "should have been recording"?
- How does the system know when a user liked/did not like an automated decision?
- Hypothesis: this is why humans have dialogue ... as in Habermas: "Theory of Communicative Action" ... but this is a tough challenge ... until then ...

206

205



Need to Mitigate Risks & Work With Them

Will discuss...

- Conceptual Foundations
- Then Current Human-Al Interaction Frameworks

Starting out where conceptual foundations of HCI/UB/UX/IxD typically trail off (after basics of perception, interaction, etc.)

20









212

211

ГЭ



Embodiment: Phenomenology

Away from Cartesian dualism...

Edmund Husserl	Martin Heidegger	Alfred Schutz	Maurice Merleau- Ponty	Ludwig Wittgenstein
Founder of Phenomenology From abstract Galilean science to things that matter Questions of experience, memory, mind, cognition	 Ending Cartesianism (separation of inner mental life and outside world) Dasein (being-in- the-world) No theory prior to praxis 	The Phenomenology of the Social World Life-world & intersubjectivity	The Phenomenology of Perception Body + phisical & social skills Embodiment	 Semiotician ,The meaning of a word is how we useit." There is no truth language games
2021-12-28		Where The Action Is		214











Reality-based interaction (RBI)



Jacob, R. J. K., Girouard, A., Hirshfield, L. M., Horn, M. S., Shaer, O., Solovey, E. T., & Zigelbaum, J. (2008). Reality-based interaction: A framework for post-WIMP interfaces. *Proceedings of the Twenty-Sixth Annual SIGCHI Conference on Human Factors in Computing Systems*, 201–210.

_



Social Context • other people • status (bosses don't type)

- showing off
- competition
- fear of failure
- motivation
 - fear
 - allegiance
 - ambition
 - self satisfaction

222

224

- inadequate systems
- frustration
- lack of positive motivation





Collaboration

- A "reason" for us being social
- Shared intentions
- Joint-action (deeply engrained)
- CSCW
- Key concepts: trust, reliability, responsibility
- We apply principles to non-human sentient beings (e.g. history of human <-> dog collaboration)

223



- Modern computers can be really good at it
- XX chromosome carriers are worse
- XY chromosome carriers are even worse!
- Everyone is actually "somewhere in between" and despite population averages individuals can be anything...

Cognitive Dimensions / Aspects

• learnings from what we know about cognition applied to design considerations...





226

228

225



Setting the Scene

- How is the system presented relative to its capabilities?
- How does user control work, how is it communicated?
- For complex interactions: are there <u>repair strategies</u>?
 Undo / "go back" etc.?
- `Expression of uncertainties?
- Living with imperfections, recovering from errors?

... can use the concepts from above to inform design thinking ...





HelpMe Robot Journeys 🛏

- hitchBOT (2013 2015)
- Can robots trust humans?
- Gained international attention for successfully hitchhiking across Canada, Germany and the Netherlands ... ☺
- ... but in 2015 its attempt to hitchhike across the United States ended prematurely when the robot was stripped and decapitated in Philadelphia, Pennsylvania (8)
- Now actively used in delivery robots!

Wood, L. J., Zaraki, A., Robins, B., & Dautenhahn, K. (2019). Developing Kaspar: A Humanoid Robot for Children with Autism. *International Journal of Social Robotics*. <u>https://doi.org/10.1007/s12369-019-00563-6</u>

Otherware

- People perceive autonomous systems / agents as counterparts
- From embodied relationship to alterity
- Technology becomes other...
- We call this class of interactive systems otherware (<u>https://otherware.net</u>)

Hassenzahl, M., Borchers, J., Boll, S., Pütten, A. R. der, & Wulf, V. (2021). Otherware: How to best interact with autonomous systems. *Interactions*, *28*(1), 54–57. https://doi.org/10.1145/3436942

232



General AI/MI Model Construction





Human-Al Interaction Design Frameworks

Big corp. research arms chiming in...

- Microsoft: <u>https://www.microsoft.com/en-</u> us/research/project/guidelines-for-human-ai-interaction/
- Google: https://pair.withgoogle.com/guidebook/
 - Focus on responsible AI ... more later
- IBM: <u>https://developer.ibm.com/technologies/machine-learning/articles/machine-learning-and-bias/</u>
 - Focus on preventing bias ...

237





238











Learning Goals

Brief outline of relevant concepts with pointers to relevant further materials:

- Contextualise motivation / need for explainable AI
- Development of the area and key principles
- Explainable / x-able Al
- Links to concepts from interaction design (e.g. UB/UX)
- Links to challenges in Al

245



246

245

247

Explainable AI (XAI)

- EU General Data Protection Regulation (GDPR) stipulates right to obtain "meaningful information about the logic involved" commonly interpreted as a "right to an explanation" for consumers affected by an automatic decision (Parliament and Council of the European Union, 2016)
- Developing field: no clear agreement about what an explanation is, nor what a good explanation entails
- First emerged in mid-1980s among wave of "expert systems" (remember: move to more practical applications)

Confalonieri, R., Coba, L., Wagner, B., & Besold, T. R. (2021). A historical perspective of explainable Artificial Intelligence. WIREs Data Mining and Knowledge Discovery, 11(1), e1391. https://doi.org/10.1002/widm.1391 Applied examples (code):

https://github.com/jphall663/interpretable_machine_learning_with_python

recommenders / neuro-symbolic learning and reasoning

Confalonieri, R., Coba, L., Wagner, B., & Besold, T. R. (2021). A historical perspective of explainable

• Different approaches in expert systems / ML /

Artificial Intelligence. WIRE bata Mining and Knowledge Discovery, 11(1), e1391. https://doi.org/10.1002/widm.1391

• Social motivation, but also commercial motivation

Explainable AI (XAI)

• Explainability can be an excuse...















Interpretability

- Explainable (frequently): blackbox with posthoc explanations
- Interpretability (frequently): model that is not a black box
- Interpretability (sometimes): Focus on outputs
- x-able Al:
 - Relatable

Good explanation/bad explanation?



Good explanation/bad explanation?



257











Exploitability

- Abusive AI use / relationships
- Not talking about "hurting agents" (yet!?)...
- Embedding models can leak information
- Developer API access can be enough to be privacy critical
 - Esp. with models trained on private data
 - E.g. Thieves of Sesame Street: Model Extraction on BERT-based APIs
 - Extracted models quickly show high correlations with orig. mod.

The Case for Usable Al

- Modern AI in video games (as in many other application areas) usually far from state of the art
- Problem solving capacity vs. usability
- Critical "real-world use" considerations:
 - Plausibility / Believability
 - Computational Performance
 - Ease of Implementation

https://www.youtube.com/watch?v=aqXHsBrS6_U

Pfau, J., Smeddinck, J. D., & Malaka, R. (2020). The Case for Usable AI: What Industry Professionals Make of Academic AI in Video Games. In *Extended Abstracts of the 2020* Annual Symposium on Computer-Human Interaction in Play (pp. 330–334). Association for Computing Machinery. <u>https://doi.org/10.1145/3383668.3419905</u>

264







Accuracy of recidivism prediction CMPAS tool (137 features): $65\% \pm 1\%$ (slightly beter (ban random) Age and number of priors Material to the transport of priors Age and number of priors Device to the transport of t

Complex Related Concepts

- Filter Bubbles and the Attention Economy
- Misinformation, Disinformation, Gossip and Bullshit ... Consider: The Social Dilemma
- Not all Al-based, but can arguably intensify feedbackloops

Al-Solutionism

- Harms of AI for predicting social outcomes
- Hunger for personal data
- Massive transfer of power from domain experts & workers to unaccountable tech companies
- Lack of explainability
- Distracts from interventions
- Veneer of accuracy

Via: Narayanan, A. (n.d.). How to recognize AI snake oil. 21.

270







Ethics History (Human Subject Research)

Mengele's twin "studies" (1940s)





Image L: https://commons.wikimedia.org/wiki/Fle:Child_survivors_of_Auschwitz.jpeg (public doma Image R: https://en.wikipedia.org/wiki/Fle:Tukegee study.jpg (public domain)

People run these things \ldots often "following orders" \ldots troubling relation to locus of responsibility / justification practices w. Al

273

General Research Ethics

Following the European Code of Conduct for Research Integrity by the ESF (as of 2011)...

- honesty in communication;
- reliability in performing research;
- objectivity;
- impartiality and independence;
- openness and accessibility;
- duty of care;
- fairness in providing references and giving credit; and
- responsibility for the scientists and researchers of the future.

274

Ethical Key Concepts

- Act utilitarianism: A person's act is morally right if and only if it produces the best possible results in that specific situation
- Rule utilitarianism: A person's act is morally right if and only if it conforms to a *rule* that leads to the greatest good.
- Deontology: The morality of an action should be based on whether that action itself is right or wrong under a series of rules, rather than based on the consequences of the action.

Ethical Key Concepts

- Virtue ethics:
 - An act is moral if it is virtuous
 - A virtue is 'indeed a character trait—that is, a disposition which is well entrenched in its possessor, something that, as we say "goes all the way down"... ' It requires the *practically wise* agent.
 - the practically wise agent's [has] capacity to recognise some features of a situation as more important than others, or indeed, in that situation, as the only relevant ones
 - 2. Practical wisdom can only come with experience (e.g. to know the likely consequences of certain actions for people)

→ Current systems can't be blamed? But they will be? Account ability? Nearly any ethical frameworks disallows "I followed therules" as defense.

Who is to blame?

"in March 2018, a self-driving Uber was navigating the Phoenix suburbs and failed to "see" a woman, hitting and killing her... In the case of Uber, the person minding the autonomous vehicle was ultimately blamed, even though Uber had explicitly disabled the vehicle's system for automatically applying brakes in dangerous situations."

via https://haiicmu.github.io from https://ainowinstitute.org/Al_Now_2018_Report.pdf

277



278







Further Considerations

- Al & Data Ethics Inseparable
- Complex relationships between ethics (focus on moral), policy & legal concerns, and commercial concerns
- Global-Local Problem
 - Diverse teams, weirdly spread out, global apps



285

Big Tech Frameworks

E.g.

- https://ai.google/principles
- <u>https://www.microsoft.com/en-us/ai/responsible-ai?activetab=pivot1%3aprimaryr6</u>
 - Fairness, Inclusiveness, Reliability & Safety, Transparency, Privacy & Security, Accountability

Review framework: <u>https://arxiv.org/abs/2001.00973</u> Closing the AI Accountability Gap: Defining an End-to-End Framework for Internal Algorithmic Auditing

Model Cards for Model Reporting: Model Cards

286


$\mathsf{Human} \rightarrow \leftarrow \mathsf{AI}$

- Humans can be (more or less) active contributors, not "just consumers"
- Humans can be involved (beyond being researchers and developers) in any step of the AI lifecycle process...
 - E.g. collect data, label examples, outcome validation, etc.
- Key concepts: crowdsourcing & human computation



"Crowd" + "Outsourcing"

- One specific area of application that leverages wisdom (or, at the very least, earnest contributions) of crowds
- Originally used to describe how businesses use the Internet to "outsource work to the crowd"
- Now also used to describe how many different types of online projects (including academic and non-profit) use crowds for social good.

(Robin Wienschendorf, 2017)

290

289

289



Motivation EXTRINSIC INTRINSIC • Fun (e.g. games) • Money Reputation • Altruism • Awards Social comparison / compet. Networking Learning PIXEL WORS Motivierende Effekte verschiedener Level der Gamification im Bereich Human-based Computation für Image Labeling und Segmentierung

Human Computation Games: Robots & Pancakes



https://youtu.be/REEgBzvcjMQ

Walther-Franks, B., Smeddinck, J., Szmidt, P., Haidu, A., Beetz, M., & Malaka, R. (2015). Robots, Pancakes, and Computer Games: Designing Serious Games for Robot Imitation Learning. Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems, 3623–3632. https://doi.org/10.1145/2702123.2702552

293

293



Taxonomy Human Computation • Human computation: means of solving computational problems • Rarely found in Social Crowdsourcing Computing crowdsourcing / social computing applications Collective Intelligence Quinn, A. J., & Bederson, B. B. (2011). Human computation: A survey and taxonomy of a growing field. *Proceedings of the 2011 Annual* Data Mining Conference on Hum an Factors in Computing Systems, 1403-1412.

Imagine that you're setting up a crowdsourcing task for image classification: "Dog or Muffin?"

What "quality control" strategies

via https://haiicmu.github.io

296

would you use to get an accurate answer for each

image?

Quality Control

296

294

Problem Spaces for HComp

- Intuitive Decisions
 - Now being partially replaced (e.g. in perception)
 - Still holds as intuitive judgment / in social contexts
- Aesthetic Judgment
- Contextual Reasoning
- Embodiment Issues

Krause, M., & Smeddinck, J. (2011). Human Computation Games: A Survey. 19th European Signal Processing Conference, 2011, 754–758.

297

Some roles for humans

- Data feeders: when machines can't read this data, but humans can.
- Backup: Machines do the easy tasks, humans do the hard stuff
- Triagers: Humans see which particular task needs to be done by people, which one can be automated (or deciding between different automation methods)
- Appeals judge: Humans override the algorithm when it's wrong
- Worker: Human does arbitrary task under machine supervision

via https://haiicmu.github.io

298

297

Some roles for machines

- Advisor: suggest (but not mandate) actions for one person
- Orchestrator: suggest (but not mandate) actions for many people
- Questioner: ask humans to consider/re-think some aspects of the task
- Manager: Tell humans what to do in task
- Lifeguard: Prevent humans from making high-cost mistakes
- Worker: Do specific aspects of the task that humans direct

via https://haiicmu.github.io

299

Humans in the Loop with AI and ML

- From data generation / curation to "live function calls"
- In ML e.g.: active learning
- Example for automated exercise execution renderings based on text:
 Sarma, H., Baran Samaddar, A., Porzel, R., D. Smeddinck, J., & Malaka, R. (2017). Updating Bayesian networks using crowds. Neural Network World, 27, 529–540. https://doi.org/10.14311/NNW.2017.27.028
- Example from self-driving cars: https://drive.google.com/file/d/1htnA4_bUdfXdpDo033Tz0 NLjMGG_wfuK/view











Recommender Systems

- What if you see an entirely new user?
 - Avoid predicting all 0 due to optimisation terms .. Mean normalisation...
- Also content-based recommendations
- Can be formulated as supervised learning problem
 Optimisation via gradient descent

306



Decision Support Systems

• E.g. in digital health:

Zhu, N., Cao, J., Shen, K., Chen, X., & Zhu, S. (2020). A Decision Support System with Intelligent Recommendation for Multi-disciplinary Medical Treatment. ACM Transactions on Multimedia Computing, Communications, and Applications, 16(1s), 33:1-33:23. https://doi.org/10.1145/3352573

• Diagnostics, ICD-codes, most viable mediation(s), etc.



















Cognitive Dimensions: Communication



317

Compression & Aliasing

- Discretisation / Quantisation
- Sampling
- Aliasing (rasterization & signal undersampling / oversampling)
- Nyquist Theorem (sample at 2x frequency for which aliasing should be avoided)
- Lossless vs. Lossy (perceptive)
- More detailed (but still quick) explanation: <u>https://www.youtube.com/watch?v=yWqrx08UeUs</u>

318



Dialogical Service Provisioning

From ongoing work with Viana (Nijia) Zhang (supervisee at Open Lab):

- Maternal mental health and well-being
- Un-platforming approach / chat-based → Focus: NLI aspects



N.Zhang10@newcastle.ac.uk

Conv. Interfaces $\leftarrow \rightarrow$ Conv. Agents

- Language != conversation
- Dialogue (social rules implied)
- Arguably building block of cognition
 See Wittgenstein: Language Games
- \rightarrow (some degree of) antropomorphologisation

See e.g.: Wang, Q., Saha, K., Gregori, E., Joyner, D. A., & Goel, A. K. (2021). Towards Mutual Theory of Mind in Human-AI Interaction: How Language Reflects What Students Perceive About a Virtual Teaching Assistant. 15. <u>http://qiaosiwang.me/Publications/MToM_Preprint.pdf</u>

322



Modern AI for NLP/NLI

Language unfolds over time \ldots need to allow methods to capture that \ldots

- Recurrent Neural Network (RNNs)
- Long Short Term Memory RNN (LSTM)
- Transformers

See also:

https://towardsdatascience.com/recurrent-neural-networks-deeplearning-for-nlp-37baa188aef5 And: https://medium.com/analytics-vidhya/natural-languageprocessing-from-basics-to-using-rnn-and-lstm-ef6779e4ae66

323



LSTM: Long Short Term Memory RNN

LSTM adds 3 gates to manage memory





The Gate of Forgetting What past hidden state is worth keeping (still relevant) or no: (0 forget - 1 keep) * past via https://hailcmu.github.io

The Gate of Input What past hidden state will be useful to figure out this input? Input + (remembered past)

The Gate of Output What out of past state + this current input will be useful later?

325



LSTM model (a neuron is now called a cell)

LSTM also adds a long term memory (cell state) to complement a more short term memory (hidden state).



326



328

via https://haiicmu.github.io

Memory and attention

What are some things this voice agent (HAL) is able to do that our voice agents are not able to do today?



via https://haiicmu.github.io

329

In recalling sequences, what information do you pay attention to?

330

332



Neural Nets with Attention



Like in this heatmap visualization, "Attention" is a technique that can be used to see what words in a sequence the model pays the most attention to make its prediction.

Shown is attention for translating a phrase from english to french



<section-header><image><text><text><text><text><text>

Crossing Domains

- Transformers
- originally developed for language problems
- Growing body of work on other domains:
 - Images
 - <u>Videos</u>
 - <u>Speech</u>
 - protein folding
 - automated coding (sort of language)
- ...



Learning Goals

Broad understanding of interaction concerns around:

- Bots & (interface) agents
- Embodied conversational agents & NPCs
- Robots & human-robot interaction



338

Bots

- Most bots are simply apps that run scripts on the internet (aka web robots aka Internet bots)
 - E.g. search engine spiders (interaction: robots.txt)
 - Many are "bad" ... e.g. "spambots"
- Often "chatbots" (i.e. conversational agents)
 - Early forms in IRC

In Prote by texysue is a locased under <u>CCRY</u>

339

Embodied Conversational Agents (ECAs): OLGA (1998)

- Chatbots with (virtual) "bodies"
- May / may not have voice
- May / may not use text





https://www.youtube.com/watch?v=Ft9oKtosMig

340

It looks like you're

Would you like help?

letter without

Don't show me

this tip again

342

help

 Get help with writing the letter
 Just type the

writing a letter.

ECAs in MR: Welbo (2000)Df:::ttp://doi.org/10.1145/633292.633299ChildrenChildrenStateDf:::ttp://www.cs.und.edu/hcil/chivideosi...Df:::ttp://www.cs.und.edu/hcil/chivideosi...State</

What killed early ECAs?

- E.g. MS Office Assistant used Bayesian Nets tech to offer support ... yay!?
- Most generally anthromorphopologisation
 - I.e. "it looks and acts at if it might be smart, so it better be smart"
 - Projected abilities often misconceptions
 - Complex rules of social behaviour / etiquette
 - Complex "rules" of dialogue
 - Some people actually liked Clippit (et al.)

Swartz, L. (2003). Why People Hate the Paperclip: Labels, Appearance, Behavior, and Social Responses to User Interface Agents. https://doi.org/10.13140/RG.2.1.2508.1047

342



- Mid 90s to early 00s wave ... but persistent issues
- Resurgence since 2010+
- Never went away, esp. through "cousin" development of non-player characters (NPCs) in video games

Dale, R. (2016). The return of the chatbots. Natural Language Engineering, 22(5), 811–817. https://doi.org/10.1017/\$1351324916000243



igure 1. Screenshot of a QuickWoZ sample scene with quick-play buttons on the top right.

Smeddinck, J., Wajda, K., Naveed, A., Touma, L., Chen, Y., Hasan, M. A., Latif, M. W., & Porzel, R. (2010). OutoKWoZ: a Multipurpose Wizardo-10-2r framework for Experiments with Embodied Conversational Agents. *Proceeding of the 14thInternational Conference on Int eligent User Interfaces* 427–428. http://doi.acm.org/10.1145/1719970.1720055

















Interaction \rightarrow Coordination \rightarrow Collaboration



Baxter https://www.youtube.com/watch?v=uLTBejGnxdE

Evolutionary Linguistics & Embodied Language Games



Luc Steels language development in robots: https://www.youtube.com/watch?v=Qh2yI-AL1V8

352

Robots in (e.g.) Health Applications



Wood, L. J., Zaraki, A., Robins, B., & Dautenhahn, K. (2019). Developing Kaspar: A Humanoid Robot for Children with Autism. International Journal of Social Robot ics. <u>https://doi.org/10.1007/s12369-019-00563-6</u>

Social robots in elderly care: https://www.youtube.com/watch?v=ppPLDEi82lg



Fard, M. J., Ameri, S., Ellis, R. D., Chinnam, R. B., Pandya, A. K., & Klein, M. D. (2018). Automated robot-assited surgical-skill evaluation: Predictal wanalytics approach. The International Journal of Medical Robot ics and Com put er Assisted Surgers 14(1), e1880. https://doi.org/10.1002/rcs.1850_353



354

353



https://www.youtube.com/watch?v=fn3KWM1kuAw

See also: Spot with an arm: https://www.youtube.com/watch?v=WvTdNwyADZc%2F&t=1120









Mixed Reality

 ... and now it is getting complicated ... <u>https://www.youtube.com/watch?v=YJg02ivYzSs</u>



Direct Manipulation \rightarrow Mind Control

- Vibrovests: remote control your dog: <u>https://www.youtube.com/watch?v=ofNjevpHdX0</u>
- EMS Remote control friends
- https://www.youtube.com/watch?v=JSfnm_HoUv4
- BCIs: for telepathy and (mediated) telekinesis <u>https://www.youtube.com/watch?v=kR1wvi2EFxA</u>
- TMS: remotely...
- Mind reading: https://www.youtube.com/watch?v=IUg-t609byg

361



- CCS-HCI
- Homo sapiens BRANCH with homo optimus?
- How will we communicate / translate?



362

361

Curiosities \rightarrow Pets \rightarrow Workers \rightarrow Collaborators

- Nearly there ...
- ... and then what?
- Which robot will be the "iPhone of robotics"?
 - Application area?
 - Capabilities?
 - Interaction modalities?











Human Deep Fakes (2018)

[Warning: video contains harsh political satire]

This video helped set off a bit of a public scare around deep fakes.



https://www.youtube.com/watch?app=desktop&v=cQ54GDm1el

via https://haiicmu.github.io

Style transfer

Neural Net is trained to learn a particular style.

We use transfer learning to apply that learned style to a new image to stylize it.



More examples: <u>https://deepart.io/latest/</u>

via https://haiicmu.github.io

369



CNN: Convolutional Neural Net

 $\ensuremath{\mathsf{CNNs}}$ are a type of multilayer perceptron most commonly used for image tasks.

We add a concept of "memory" to process serial data. What kind of abilities might a neural net need for vision?





Croatian sheepdog via https://haiicmu.github.io

Costa Rica stray dog estate

Herding pup in training 370

370









CNN building blocks: Pooling

Pooling is the opposite to padding.

A pooling layer takes a convoluted feature after the convolution layer, and shrinks it again.

Max pooling = just take the max value, throw out the rest.

Purpose: reduce computational complexity and prevent overfitting.



via https://haiicmu.github.io







Generative Adversarial Network: generating stuff

Fake paintings: Photo via Art and Artificial Intelligence Laboratory, Rutgers University

























Al & Data Visualisation



Al & (Creative) Language







394

"AI"



Human-Al Co-Creation (development) **New Project** Open an existing project from Drive. Open an existing project from a file. Image Project **Audio Project Pose Project** Teach based on images, from Teach based on one-second-long Teach based on images, from sounds, from files or your files or your webcam. files or your webcam. microphone. https://teachablemachine.withgoogle.com/train 396







Course Outline: Week 02

- Topic 01 Machine Learning Basics
- Topic 02 Traditional Machine Learning Models
- Topic 03 Machine/Deep Learning for Human Activity Recognition.
- Topic 04 Research on Wearable-based Behaviour Analysis

Course Outline: Week 03

- Topic 01 Explainable, Interpretable & Relatable AI
- Topic 02 AI Ethics
- Topic 03 Humans-in-the-Loop
- Topic 04 Recommender Systems
- Topic 05 Conversational Interfaces
- Topic 06 Al Agents & Robots
- Topic 07 Human-Al Integration
- Topic 08 Creative Al
- Topic 09 Summary & Outlook

401



401













Mental Models

- **<u>O Set expectations for adaptation</u>**
 - Identify and building on existing mental models
 - "What is the user trying to do?"
 - "What mental models might already be in place?"
 - "Does this product break intuitive patterns?"

• **Onboard in stages**

- Set realistic expectations early
- Describe user benefits, not technology
- Describe the core value initially, introduce new features as they are used
- Make it easy for users to experiment

• **③** Plan for co-learning

- Connect feedback to personalization and adaptation
- Fail gracefully to non-Al options when needed
- O Account for user expectations of human-like interaction
 - Communicate algorithmic nature and limits of system

409

Explainability + Trust

• **<u>O Help users calibrate their trust</u>**

- Instill user trust it in some situations, but indicate to double-check when needed
- Articulate data sources
- Tie explanations to user actions
- Optimize for understanding
 - Consider partial / full explanations, and/or process insights
- 3 Manage influence on user decisions
 - Consider communicating model confidence

410

409

411

Feedback + Control

- **O** Align feedback with model improvement
 - Clarify differences between implicit and explicit feedback
 - Asking the right questions at the right level of detail

• 2 Communicate value & time to impact

- Understanding why people give feedback
- Build on existing mental models to explain benefits
- Communicate how/when user feedback will change experience

• ③ Balance control & automation

- Help users control the aspects of the experience they want to
- Easy opting out of giving feedback

Errors + Graceful Failure

- ① Define "errors" & "failure"
 - System working as intended can still be perceived as failure
 - Acknowledge work-in-progress (if applicable)

• **O Identify error sources**

- Inherent complexity can make identifying source of errors challenging
- Consider error / error source discovery strategies

• ③ Provide paths forward from failure

- Probabilistic systems will fail at some point
- Identify failure; user-centered discovery / resolve
- Facilitate feedback and return to task







An Apt Closing Statement

